

Goal: Smoother Density Estimation.

3 Kernel Density Estimation

Recall the definition of a density function

$$f(y) \equiv \frac{d}{dy} F(y) \equiv \lim_{h \rightarrow 0} \frac{F(y+h) - F(y-h)}{2h} = \lim_{h \rightarrow 0} \frac{F(y+h) - F(y)}{h},$$

where $F(x)$ is the cdf of the random variable Y .

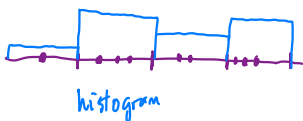
small
take this and approximate w/ fixed h ,
use Ecdf.

$$\hat{f}(x) = \frac{\hat{F}_n(y+h) - \hat{F}_n(y)}{h} \quad \text{histogram.}$$

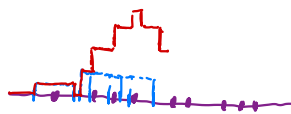
What if instead, we replace $F(y+h) - F(y-h)$ with ecdf

$$\begin{aligned} \Rightarrow \hat{f}(y) &= \frac{\hat{F}_n(y+h) - \hat{F}_n(y-h)}{2h} = \frac{\#\{y_i \in (y-h, y+h]\}}{2nh} \\ &= \frac{\sum_{i=1}^n \mathbb{I}(y_i \in (y-h, y+h])}{2nh} \\ &= \frac{1}{nh} \sum_{i=1}^n \frac{1}{2} \mathbb{I}(y-h < y_i \leq y+h) \\ &= \frac{1}{nh} \sum_{i=1}^n \frac{1}{2} \mathbb{I}\left(-1 < \frac{y-y_i}{h} \leq 1\right) \\ &= \frac{1}{nh} \sum_{i=1}^n K\left(\frac{y-y_i}{h}\right) \quad \text{where } K \text{ is a Uniform density on } [-1, 1] \end{aligned}$$

kernel function.



histogram



still not continuous (because Uniform density not continuous!)



\Rightarrow another kernel may lead to smoother estimate.

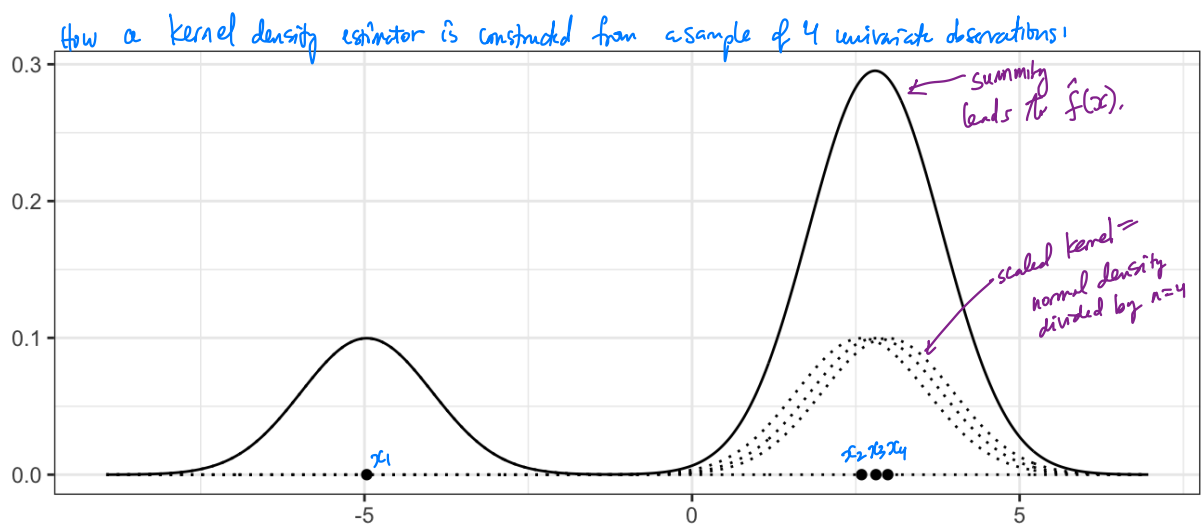
A kernel function assigns weights to the contribution given by each y_i to $\hat{f}(y)$, depending on proximity to y .

This will weight all points within h of y equally. A univariate kernel density estimator will allow a more flexible weighting scheme.

Typically, kernel functions are positive everywhere and symmetric about zero.

Examples of ideas for such functions? Normal density, Student t (others exist).

Additionally, constraining K so that $\int z^2 K(z) dz = 1$ allows h to play role of scale parameter (not required).

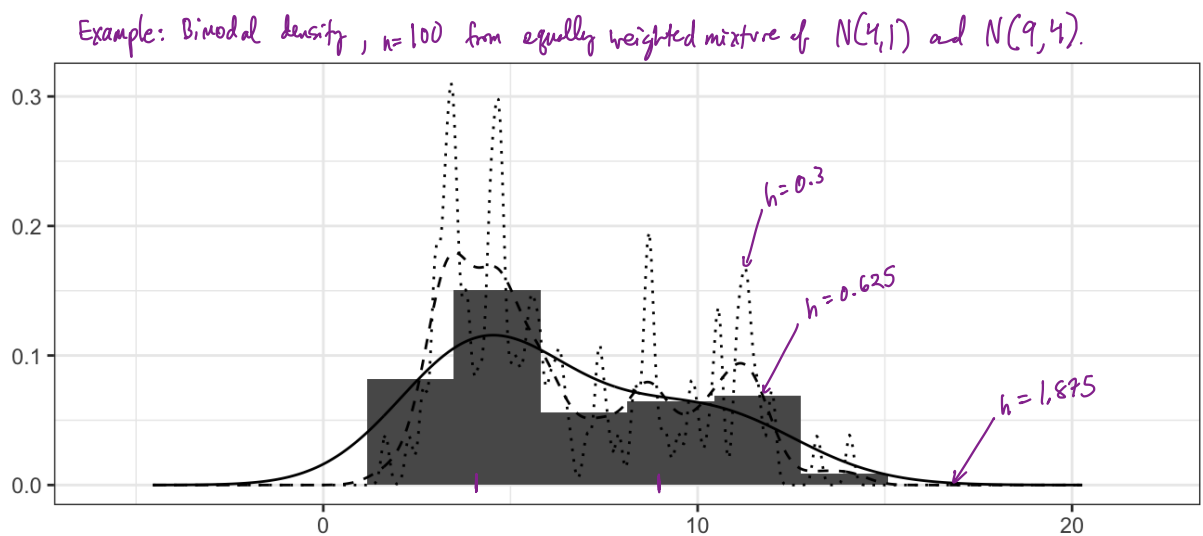


3.1 Choice of Bandwidth

The bandwidth parameter controls the smoothness of the density estimate. *for a given kernel.*

bandwidth determines tradeoff bet/ bias and variance.

The tradeoff that results from choosing the bandwidth + kernel can be quantified through a measure of accuracy of \hat{f} , such as MISE.



For large h , oversmoothing (lose 2nd mode).

For small h , undersmoothing (many false modes).

To understand bandwidth selection, let us analyze MISE ^(AMISE). Suppose that K is a symmetric, continuous probability density function with mean 0 and variance $0 < \sigma_K^2 < \infty$. Let $R(g) = \int g^2(z) dz$. Recall that

$$\text{MISE}(h) = \int \text{MSE}(\hat{f}(x)) dx = \int \text{var}\{\hat{f}_h(x)\} + [\text{bias}\{\hat{f}_h(x)\}]^2 dx$$

Now let $h \rightarrow 0$ and $nh \rightarrow \infty$ as $n \rightarrow \infty$.

Bias: Note $E\{\hat{f}_h(x)\} = \frac{1}{h} \int K\left(\frac{x-u}{h}\right) f(u) du$
 $= \int K(t) f(x-ht) dt$ (change of variable).

and using Taylor's expansion, $f(x-ht) = f(x) - ht f'(x) + \frac{h^2 t^2}{2} f''(x) + o(h^2)$

$\Rightarrow E\{\hat{f}_h(x)\} = f(x) + \frac{h^2 \sigma_K^2}{2} f''(x) + o(h^2)$ ← because K is symmetric about 0.

so, $[\text{bias}\{\hat{f}_h(x)\}]^2 = \frac{h^4 \sigma_K^4}{4} [f''(x)]^2 + o(h^4)$.

$\Rightarrow \text{IB} = \int [\text{bias}\{\hat{f}_h(x)\}]^2 dx = \frac{h^4 \sigma_K^4}{4} R(f'') + o(h^4)$.

Variance: $\text{Var}\{\hat{f}_h(x)\} = \frac{1}{n} \text{Var}\left\{\frac{1}{h} K\left(\frac{x-X_i}{h}\right)\right\}$
 $\stackrel{\text{change of variable}}{=} \frac{1}{nh} \int K(t)^2 f(x-ht) dt - \frac{1}{n} \left[E\left\{\frac{1}{h} K\left(\frac{x-X_i}{h}\right)\right\}\right]^2$
 $\stackrel{\text{Taylor's expansion}}{=} \frac{1}{nh} \int K(t)^2 [f(x) + o(h)] dt - \frac{1}{n} [f(x) + o(h)]^2$
 $= \frac{1}{nh} f(x) R(K) + o\left(\frac{1}{nh}\right)$

$\Rightarrow \text{IV} = \int \text{var}\{\hat{f}_h(x)\} dx = \frac{R(K)}{nh} + o\left(\frac{1}{nh}\right)$.

and $\text{MISE} = \underbrace{\frac{R(K)}{nh} + \frac{h^4 \sigma_K^4}{4} R(f'')}_{\text{AMISE}} + o\left(\frac{1}{nh} + h^4\right)$.

To minimize AMISE with respect to h , seek value of h that avoid excessive bias and variance.

$$\text{optimal bandwidth } h_0 = \left(\frac{R(K)}{n \sigma_K^4 R(f'')} \right)^{1/5}$$

$$\Rightarrow \text{minimal AMISE: } \text{AMISE}_0 = \frac{5}{4} \left[\sigma_K R(K) \right]^{4/5} R(f'')^{1/5} n^{-4/5}$$

$$\text{recall for histogram, } \text{AMISE}_0 = \left[\frac{R(f')}{16} \right]^{1/3} n^{-2/3}$$

\Rightarrow w kernel estimate, getting closer to parametric rate of n^{-1}

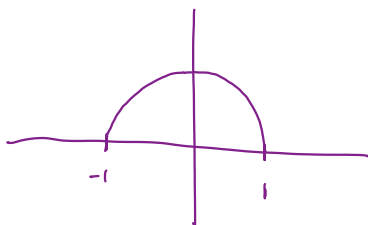
The term $R(f'')$ measures the roughness of the true underlying density. In general, rougher densities are more difficult to estimate and require smaller bandwidth.

The term $[\sigma_K R(K)]^{4/5}$ is a function of the kernel function K .

We could choose K to minimize $[\sigma_K R(K)]^{4/5}$:

If K restricted to be a proper density (w/ some moment conditions), minimizer is scaled version of a quadratic density.

$$K(u) = \frac{3}{4} (1 - u^2) \mathbb{I}(|u| \leq 1)$$



"Epanechnikov kernel" (more later).

3.1.1 Cross Validation - Want to evaluate quality of \hat{f} without using data twice (once for fitting \hat{f} , once for evaluating).

\Rightarrow again use $\hat{f}_{-i}(x_i) = \frac{1}{h(n-1)} \sum_{j \neq i} K\left(\frac{x_i - x_j}{h}\right)$ and let $\hat{Q}(h)$ be a function of $\hat{f}_{-i}(x_i)$ that assesses quality.
 \uparrow
 e.g. ISE from before.

If we choose $\hat{Q} = \text{ISE}$, choose h to minimize $R(\hat{f}) - \frac{2}{n} \sum_{i=1}^n \hat{f}_{-i}(x_i)$

\uparrow
 could instead choose $\hat{Q}(h)$ as the pseudo likelihood $PL(h) = \prod_{i=1}^n \hat{f}_{-i}(x_i)$ and choose h to maximize!

Typically under smoothed \Rightarrow too bumpy.

3.1.2 Plug-in Methods

If the reference density f is Gaussian and a Gaussian kernel K is used,

$$h_0 = 1.059 \sigma n^{-1/5}$$

\uparrow sample variance, or IQR plug-in approach.

Typically oversmoothed
 Empirical estimation of $R(f'')$ may be a better option.

$$\text{recall } h_0 = \left(\frac{R(K)}{n \sigma_K^4 R(f'')} \right)^{1/5}$$

We could use a kernel density estimator for f'' :

$$\begin{aligned} \hat{f}''(x) &= \frac{d^2}{dx^2} \left\{ \frac{1}{nh_1} \sum_{i=1}^n L\left(\frac{x-x_i}{h_1}\right) \right\} \\ &= \frac{1}{nh_1^3} \sum_{i=1}^n L''\left(\frac{x-x_i}{h_1}\right) \end{aligned}$$

where L is a sufficiently differentiable kernel and h_1 = bandwidth to estimate f' .

Note: estimating f and f'' (for $R(f'')$) will require different bandwidths.

2 stage approach [Sheather-Jones method]:

- ① Use Gaussian plug in for h_1 . This bandwidth is used to estimate $R(f'')$.
- ② h is calculated using $h = \left(\frac{R(K)}{n \sigma_K^4 R(f'')} \right)^{1/5}$ used to produce the final kernel density estimate.

3.2 Choice of Kernel

There are two choices we have to make to perform density estimation:

Kernel and bandwidth.

The shape of kernel is much less important than bandwidth.

3.2.1 Epanechnikov Kernel

Recall $AMISE_0 = \frac{5}{4} [\sigma_K R(K)]^{4/5} R(f'')^{1/5} n^{-4/5}$.

The *Epanechnikov kernel* results from choosing K to minimize $[\sigma_K R(K)]^{4/5}$, restricted to be a symmetric density with finite moments and variance equal to 1

$$\Rightarrow K(u) = \frac{3}{4} (1 - u^2) \mathbb{I}(|u| \leq 1)$$

$$\text{and } \sigma_K R(K) = \frac{3}{5\sqrt{5}}$$

\Rightarrow The ratio of $\frac{\sigma_K R(K)}{3/5\sqrt{5}}$ provides a measure of relative inefficiency of other kernels.

↓
multiplicative factor for equivalent sample size needed to achieve same AMISE.

Kernel	Inefficiency
Epanechnikov	1
Gaussian	1.0513
Uniform	1.0758
Birectangular	1.0061
Triweight	1.0135

kernel choice doesn't make much difference!

$$\frac{15}{16} (1 - u^2)^2$$

$$\frac{35}{32} (1 - u^2)^3$$

3.2.2 Canonical Kernels

Unfortunately a particular value of h corresponds to a different amount of smoothing depending on which kernel is being used.

Let h_K and h_L denote the bandwidths that minimize AMISE when using symmetric kernel densities K and L . Then,

$$\frac{h_K}{h_L} = \frac{(R(K)/\sigma_K^4)^{1/5}}{(R(L)/\sigma_L^4)^{1/5}} = \frac{\delta(K)}{\delta(L)}.$$

\Rightarrow to change from bandwidth h for kernel K to a bandwidth that gives an equivalent amount of smoothing for L , use bandwidth $h\delta(L)/\delta(K)$.

Epanechnikov: $\delta(K) \propto 15^{1/5}$, Gaussian: $\delta(K) = \left(\frac{1}{\sqrt{2\pi}}\right)^{1/5}$, Uniform: $\delta(K) = \left(\frac{9}{2}\right)^{1/5}$.

Suppose we rescale a kernel shape so that $h = 1$ corresponds to a bandwidth of $\delta(K)$,

The kernel density estimator can be written as

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n K_{h\delta(K)}(x-x_i) \quad \text{where} \quad K_{h\delta(K)}(z) = \frac{1}{h\delta(K)} K\left(\frac{z}{h\delta(K)}\right).$$

$K_{\delta(K)}$ is called a "canonical kernel"

\downarrow

benefit: a single value of h can be used interchangeably for each canonical kernel without affecting the amount of smoothing.

For a canonical kernel w/ bandwidth h ,

$$\text{AMISE} = (\sigma_K R(K))^{4/5} \left(\frac{1}{nh} + \frac{h^4 R(f'')}{4} \right)$$

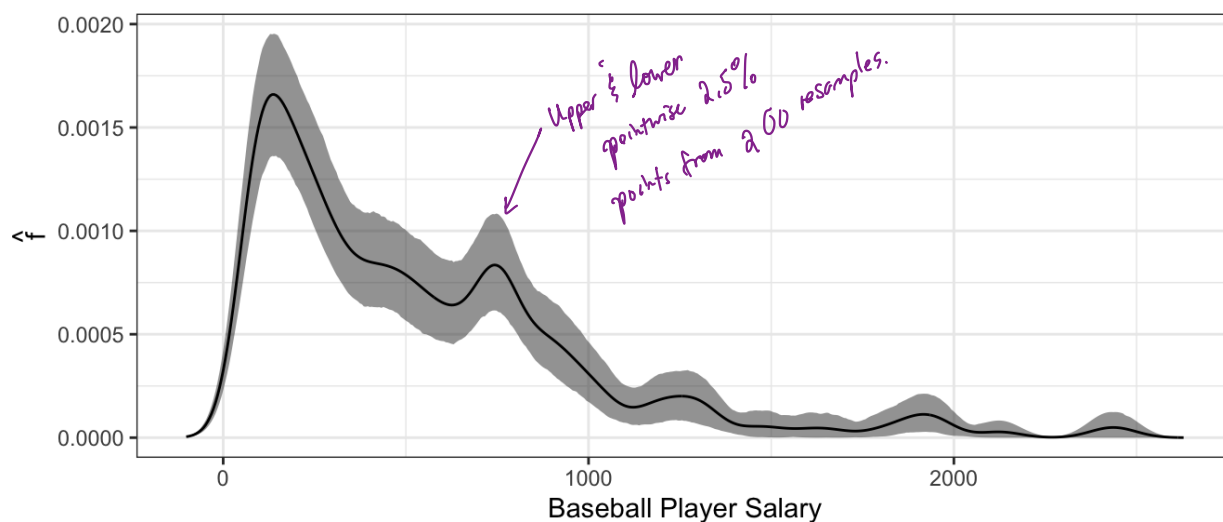
$\underbrace{\hspace{10em}}$
 bias/variance trade-off no longer confounded w/ $R(K)$ (choice of kernel).
 AND optimal kernel choice doesn't depend on bandwidth.

so, the Epanechnikov kernel shape is optimal for any desired degree of smoothing.

3.3 Bootstrapping and Variability Plot

describing uncertainty in \hat{f}

- ① A sample of size n is drawn w/ replacement from data.
- ② bandwidth chosen for new sample based on Sheather-Jones method, density estimate determined.
- ③ Repeat ①-② many times w/ values of \hat{f} recorded at a fixed grid of values.



Note: This is NOT a 95% CI for f , but a representation of the variability in the process of estimating \hat{f} .

↳ can't say anything about $P(\hat{f}(x)_{\text{lower}} \leq f(x) \leq \hat{f}(x)_{\text{upper}})$ because bias.

widens at peaks and valleys, narrows where \hat{f} more flat,

↳ consistent w/ MSE of $\hat{f}(x)$ directly related to $[f''(x)]^2$